

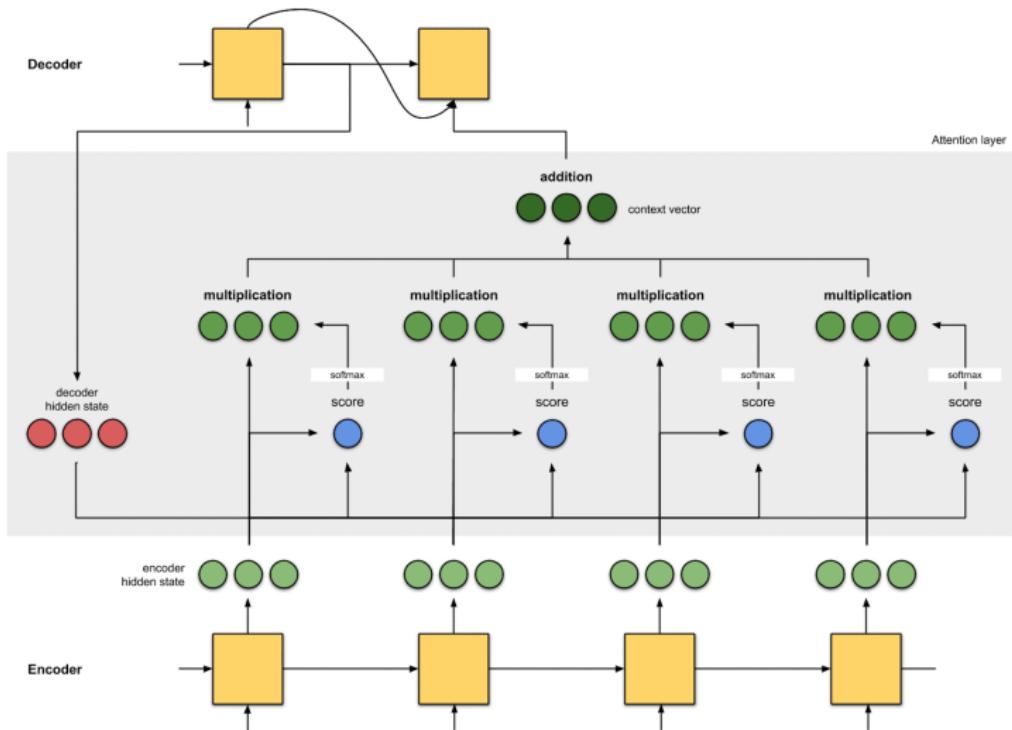
# CS391L: Machine Learning WB

## Transformers and Attention

Inderjit Dhillon  
UT Austin

April 7, 2025

# Last time: Attention + RNN in NMT



(Figure from <https://towardsdatascience.com/>)

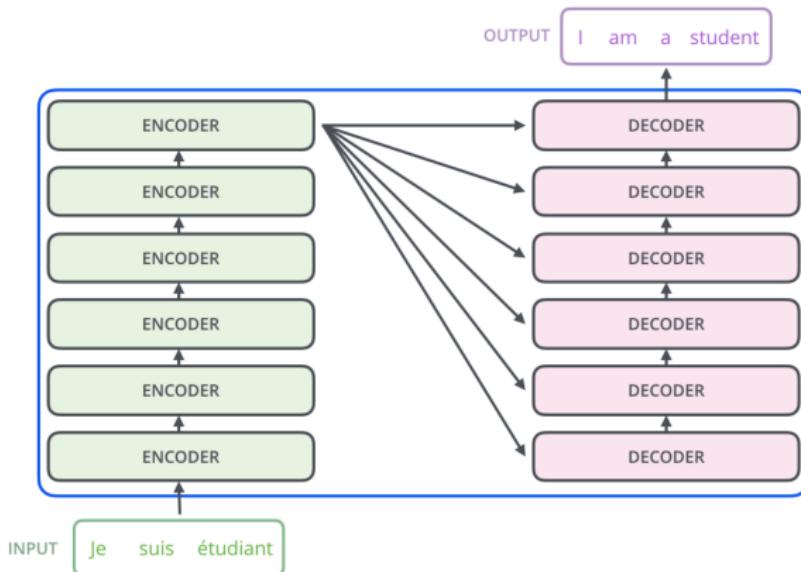
# Transformer

# Transformer

- An architecture that relies entirely on attention without using CNN/RNN
- A brief history:
  - “Attention Is All You Need” (Vaswani et al., 2017)  
First Transformer for machine translation
  - BERT (Jacob et al., 2018)  
Transformer + Pretraining achieves SOTA on many other NLP tasks.
  - Vision Transformer (Dosovitskiy et al., 2020)  
Transformer outperforms ResNet on ImageNet

# Transformer for machine translation

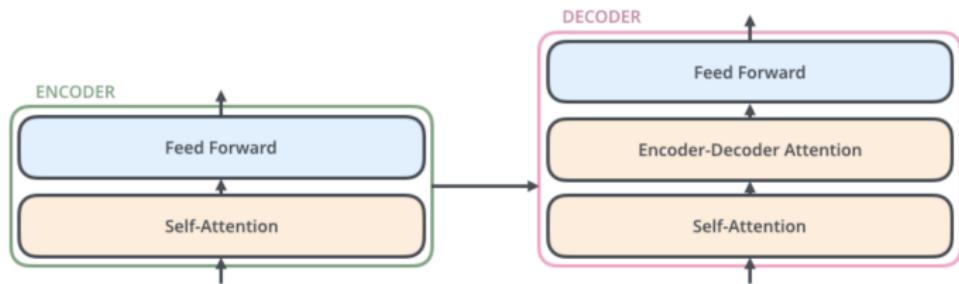
- Pass all input tokens to encoder **simultaneously**
- Passing through several Transformer blocks



(Vaswani et al., 2017)

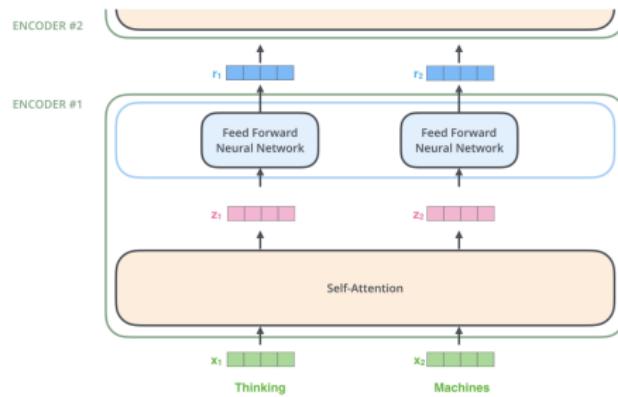
# Encoder and Decoder

- Self attention layer: the main architecture used in Transformer
- Decoder: will have another attention layer to help it focus on relevant parts of input sentences.



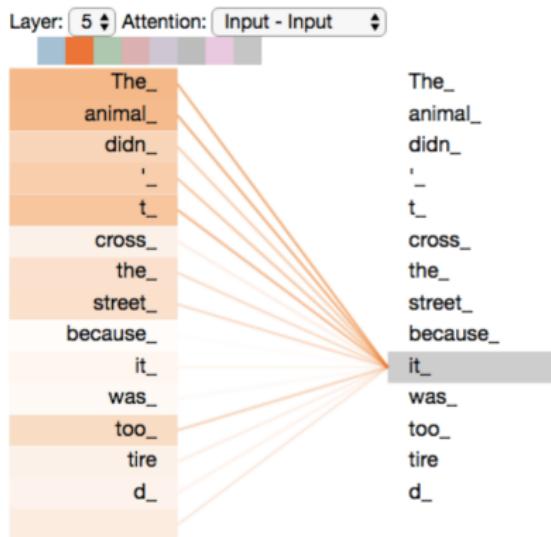
# Encoder

- Each word has a corresponding “latent vector” (initially the word embedding for each word)
- Each layer of encoder:
  - Receive a list of vectors as input
  - Passing these vectors to a **self-attention** layer
  - Then passing them into a feed-forward layer
  - Output a list of vectors



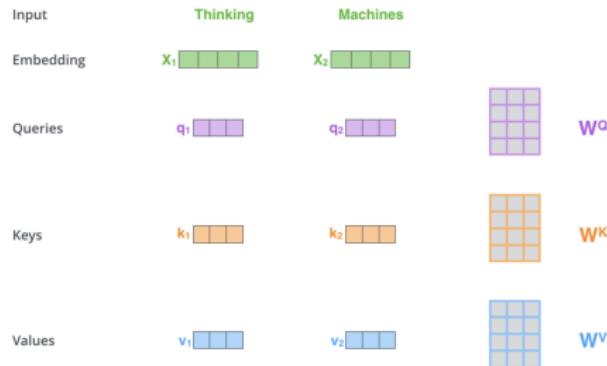
# Self-attention layer

- Main idea: The actual meaning of each word may be related to other words in the sentence
- The actual meaning (latent vector) of each word is a weighted (attention) combination of other words (latent vectors) in the sentence



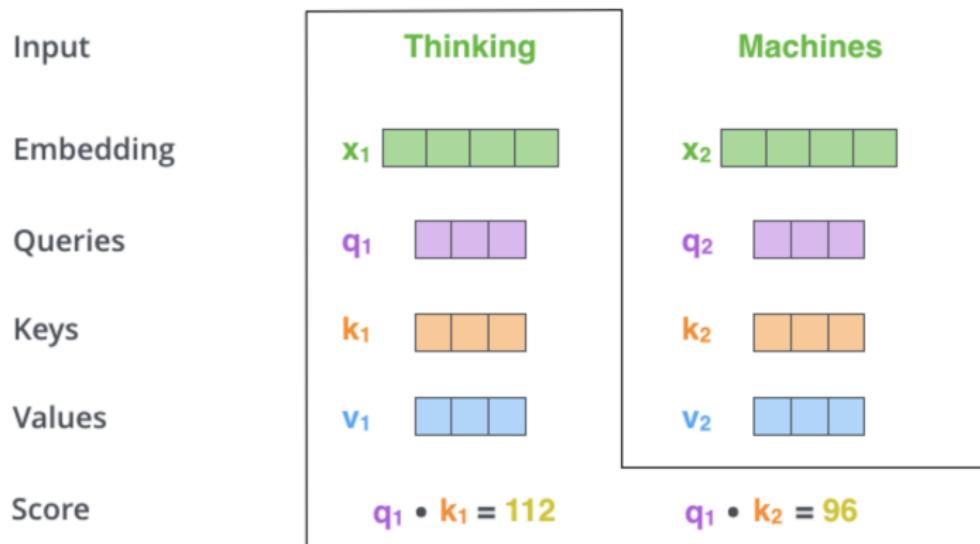
# Self-attention layer

- Input latent vectors:  $x_1, \dots, x_n$
- Self-attention parameters:  $W^Q, W^K, W^V$  (weights for query, key, value)
- For each word  $i$ , compute
  - Query vector:  $q_i = x_i W^Q$
  - Key vector:  $k_i = x_i W^K$
  - Value vector:  $v_i = x_i W^V$



# Self-attention layer

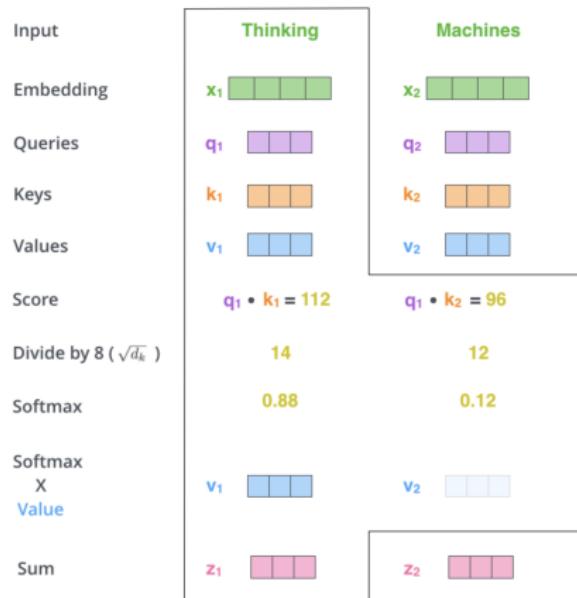
- For each word  $i$ , compute the scores to determine how much focus to place on other input words
  - The **attention** score for word  $j$  to word  $i$ :  $q_i^T k_j$



# Self-attention layer

- For each word  $i$ , the output vector

$$\sum_j s_{ij} \mathbf{v}_j, \quad \mathbf{s}_i = \text{softmax}(\mathbf{q}_i^T \mathbf{k}_1, \dots, \mathbf{q}_i^T \mathbf{k}_n)$$



# Matrix form

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V, \quad Z = \text{softmax}(QK^T)V$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^Q \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} Q \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

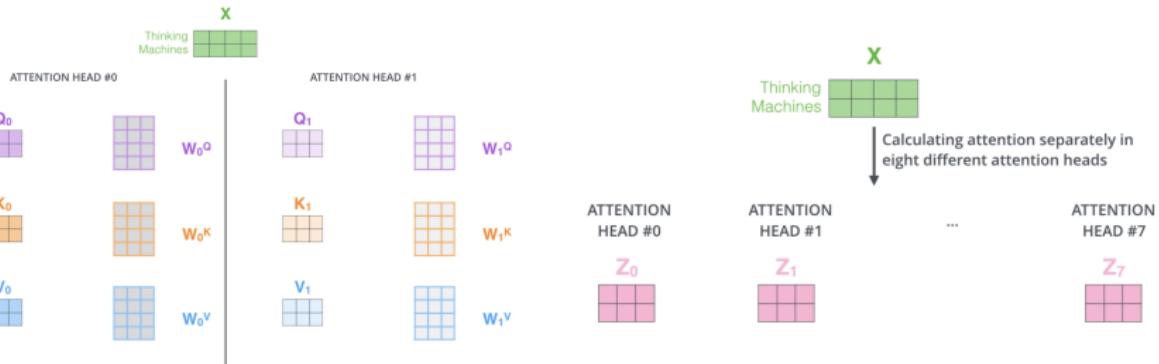
$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^K \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} K \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} W^V \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix} = \begin{matrix} V \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

$$\text{softmax}\left(\frac{\begin{matrix} Q \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix} \times \begin{matrix} K^T \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) = \begin{matrix} Z \\ \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline \end{array} \end{matrix}$$

# Multiple heads

- Multi-headed attention: use multiple set of (*key*, *value*, *query*) weights
- Each head will output a vector  $Z_i$



# Multiply with weight matrix to reshape

- Gather all the outputs  $Z_1, \dots, Z_k$
- Multiply with a weight matrix to reshape
- Then pass to the next fully connected layer

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model

$X$

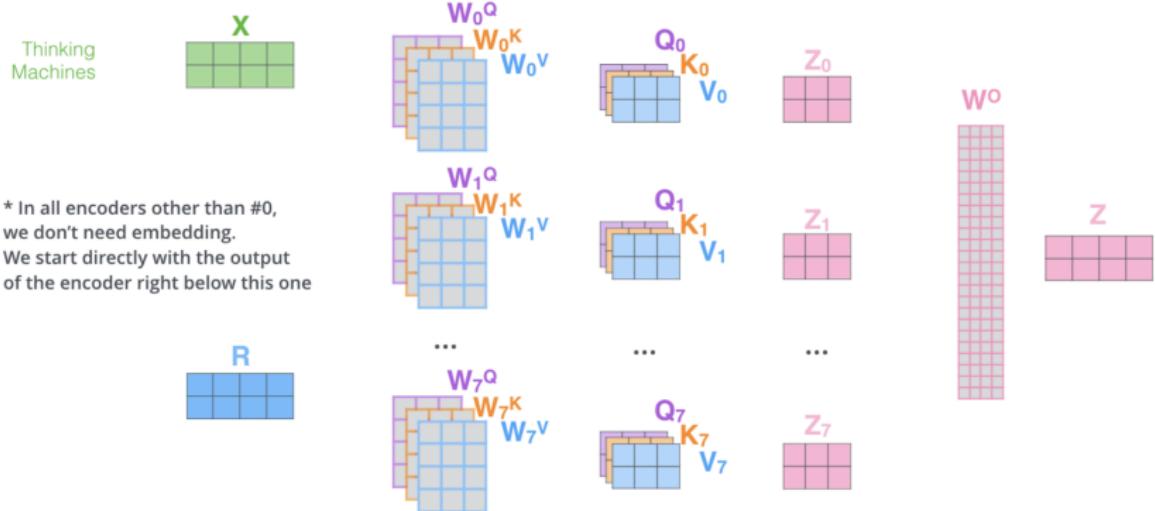


3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} \end{matrix}$$

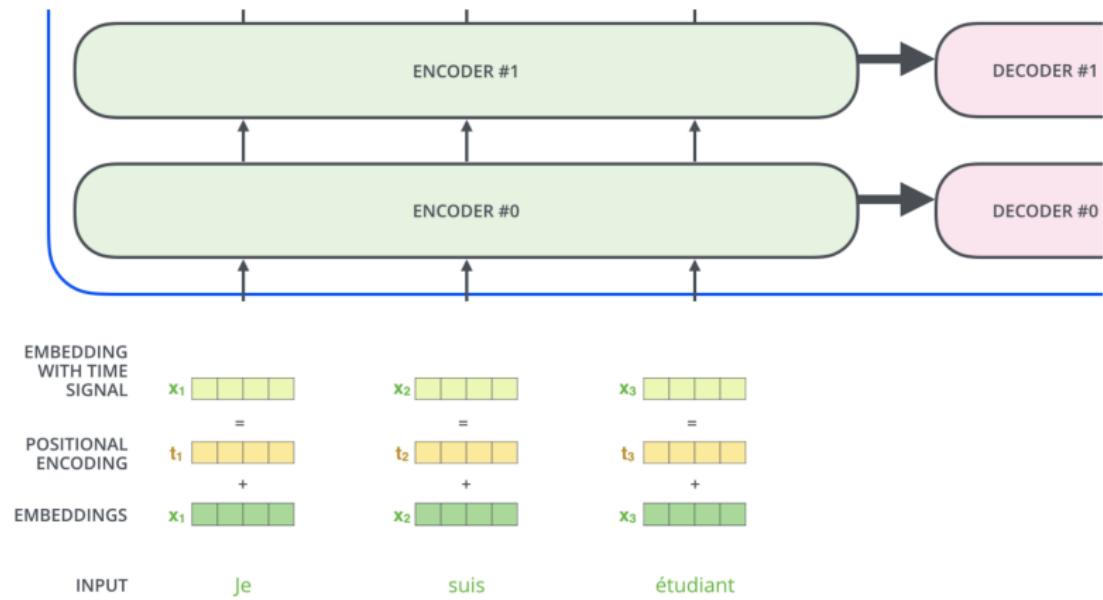
# Overall architecture

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^o$  to produce the output of the layer



# Position Encoding

- The above architecture ignores the sequential information
- Add a positional encoding vector to each  $x_i$  (according to  $i$ )



# Performance gain on NMT

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		<b><math>3.3 \cdot 10^{18}</math></b>
Transformer (big)	<b>28.4</b>	<b>41.8</b>		$2.3 \cdot 10^{19}$

(Figure from Vaswani et al., 2017)

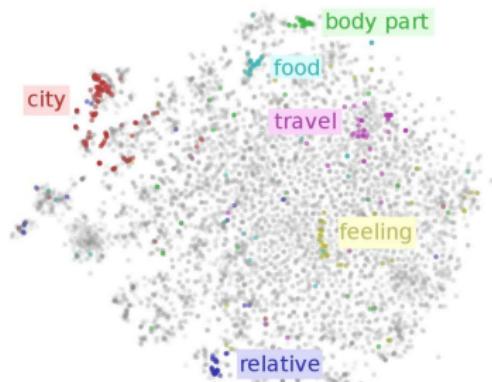
# Unsupervised pretraining for NLP

# Motivation

- There is a huge amount of unlabeled NLP data available but very little labeled data
- Can we use large amount of unlabeled data to obtain meaningful representations of words/sentences?

# Learning word embeddings

- Use large (unlabeled) corpus to learn a useful word representation
  - Learn a vector for each word based on the corpus
  - Hopefully the vector represents some semantic meaning
  - Can be used for many tasks
    - Replace the word embedding matrix for DNN models for classification/translation
- Two different perspectives but led to similar results:
  - Glove (Pennington et al., 2014)
  - Word2vec (Mikolov et al., 2013)



## Context information

- Given a large text corpus, how to learn low-dimensional features to represent a word?
- For each word  $w_i$ , define the “context” of the word as the words surrounding it in an  $L$ -sized window:

$$w_{i-L-2}, w_{i-L-1}, \underbrace{w_{i-L}, \dots, w_{i-1}}_{\text{context of } w_i}, w_i, \underbrace{w_{i+1}, \dots, w_{i+L}}_{\text{context of } w_i}, w_{i+L+1}, \dots$$

- Get a collection of (word, context) pairs, denoted by  $D$ .

# Examples

## Source Text

The quick brown fox jumps over the lazy dog. →

## Training Samples

(the, quick)  
(the, brown)

The quick brown fox jumps over the lazy dog. →

(quick, the)  
(quick, brown)  
(quick, fox)

The quick brown fox jumps over the lazy dog. →

(brown, the)  
(brown, quick)  
(brown, fox)  
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

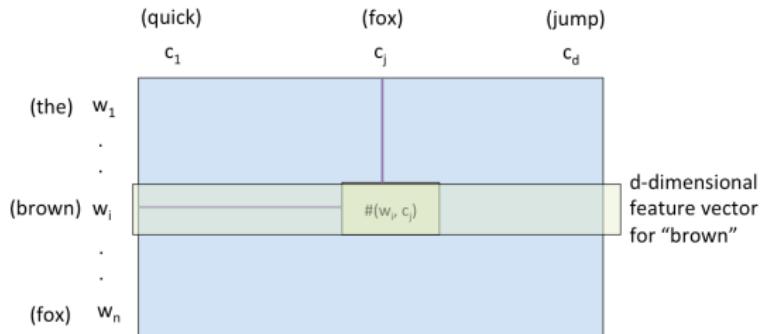
(fox, quick)  
(fox, brown)  
(fox, jumps)  
(fox, over)

(Figure from <http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>)

# Use bag-of-word model

- Idea 1: Use the bag-of-word model to “describe” each word
- Assume we have context words  $c_1, \dots, c_d$  in the corpus, compute
$$\#(w, c_i) := \text{number of times the pair } (w, c_i) \text{ appears in } D$$
- For each word  $w$ , form a  $d$ -dimensional (sparse) vector to describe  $w$

$$\#(w, c_1), \dots, \#(w, c_d),$$



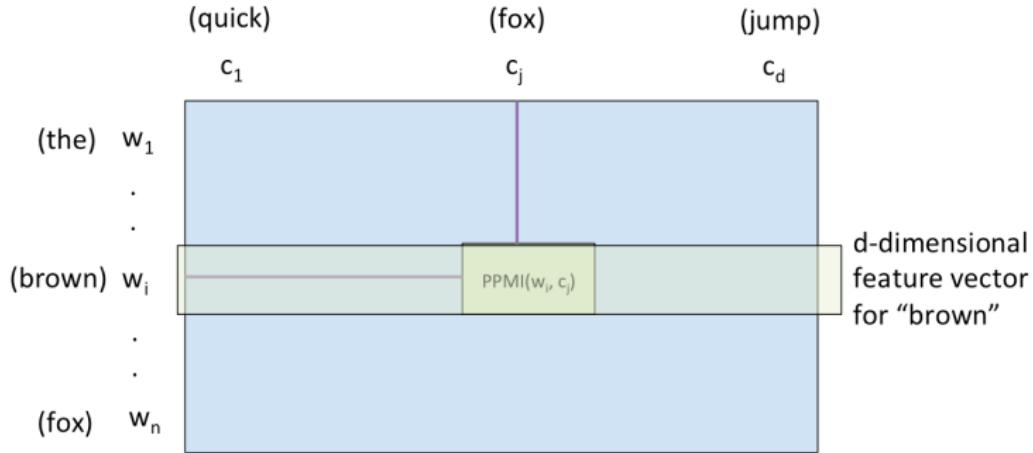
# PMI/PPMI Representation

- Similar to TF-IDF: Need to consider the frequency of each word and each context
- Instead of using co-occurrence count  $\#(w, c)$ , we can define pointwise mutual information:

$$\text{PMI}(w, c) = \log\left(\frac{\hat{P}(w, c)}{\hat{P}(w)\hat{P}(c)}\right) = \log \frac{\#(w, c)|D|}{\#(w)\#(c)},$$

- $\#(w) = \sum_c \#(w, c)$ : number of times word  $w$  occurred in  $D$
- $\#(c) = \sum_w \#(w, c)$ : number of times context  $c$  occurred in  $D$
- $|D|$ : number of pairs in  $D$
- Positive PMI (PPMI) usually achieves better performance:
$$\text{PPMI}(w, c) = \max(\text{PMI}(w, c), 0)$$
- $M^{\text{PPMI}}$ : a  $n$  by  $d$  word feature matrix, each row is a word and each column is a context

# PPMI Matrix



# Generalized Low-rank Embedding

- SVD basis will minimize

$$\min_{W,V} \|M^{\text{PPMI}} - WV^T\|_F^2$$

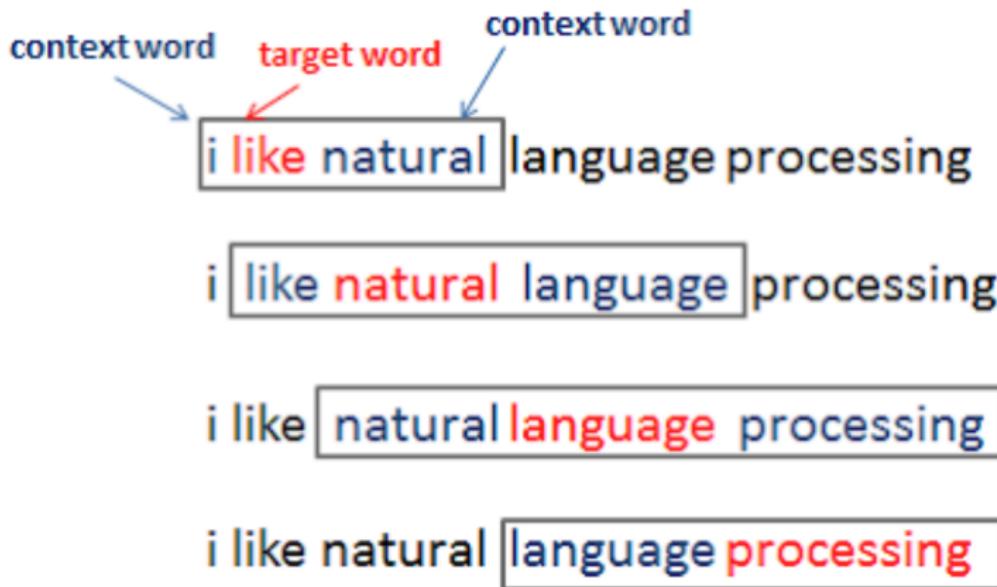
- Glove (Pennington et al., 2014)
  - Negative sampling (less weights to 0s in  $M^{\text{PPMI}}$ )
  - Adding bias term:

$$M^{\text{PPMI}} \approx WV^T + \mathbf{b}_w \mathbf{e}^T + \mathbf{e} \mathbf{b}_c^T$$

- Use  $W$  as the word embedding matrix

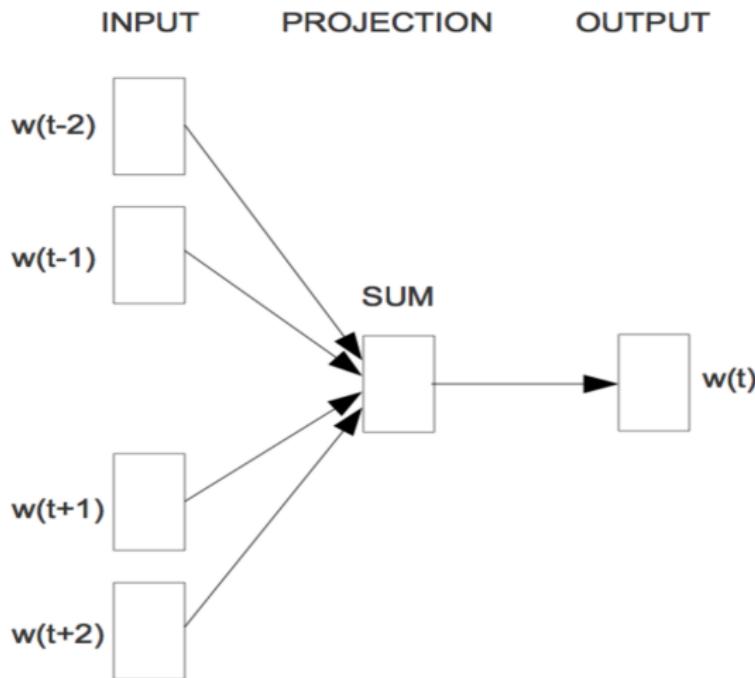
## Word2vec (Mikolov et al., 2013)

- A neural network model for learning word embeddings
- Main idea:
  - Predict the target words based on the neighbors (CBOW)
  - Predict neighbors given the target words (Skip-gram)



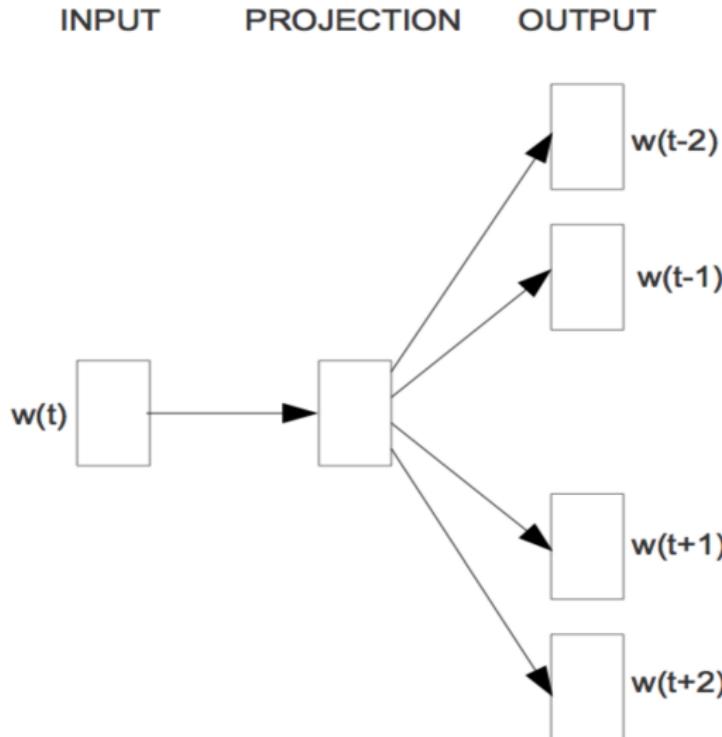
# CBOW (Continuous Bag-of-Word model)

- Predict the target words based on the neighbors



# Skip-gram

- Predict neighbors using target word



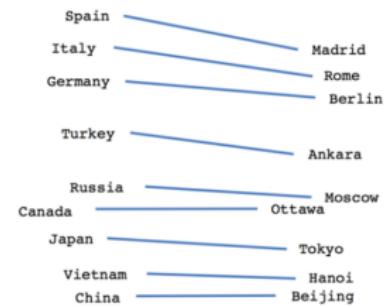
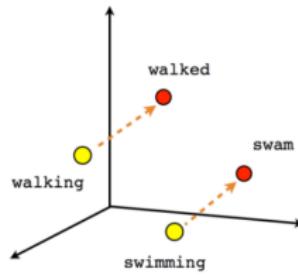
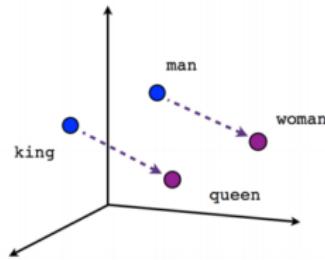
## More on skip-gram

- Learn the probability  $P(w_{t+j}|w_t)$ : the probability to see  $w_{t+j}$  in target word  $w_t$ 's neighborhood
- Every word has two embeddings:
  - $v_i$  serves as the role of target
  - $u_i$  serves as the role of context
- Model probability as softmax:

$$P(o|c) = \frac{e^{u_o^T v_c}}{\sum_{w=1}^W e^{u_w^T v_c}}$$

# Results

The low-dimensional embeddings are (often) meaningful:

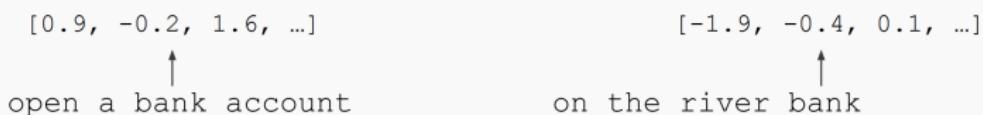


(Figure from <https://www.tensorflow.org/tutorials/word2vec>)

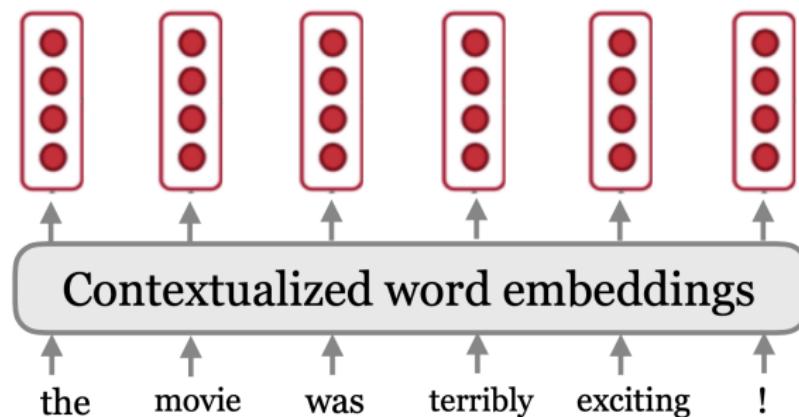
# Contextual embedding

# Contextual word representation

- The semantic meaning of a word should depend on its context

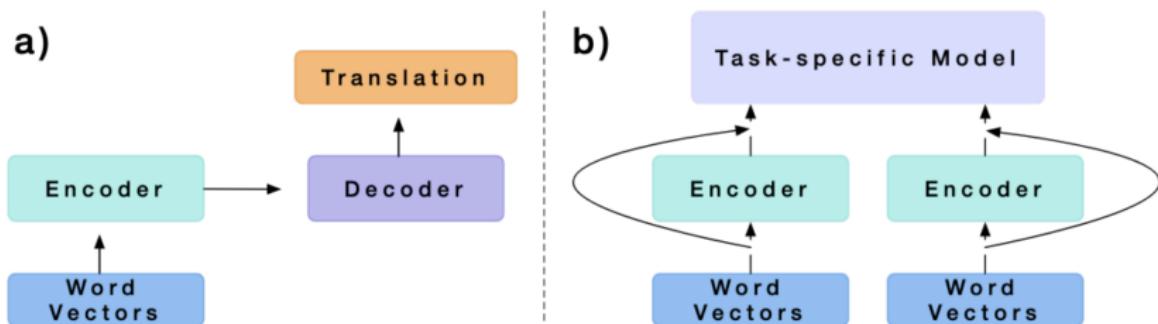


- Solution: Train a model to extract contextual representations on text corpus



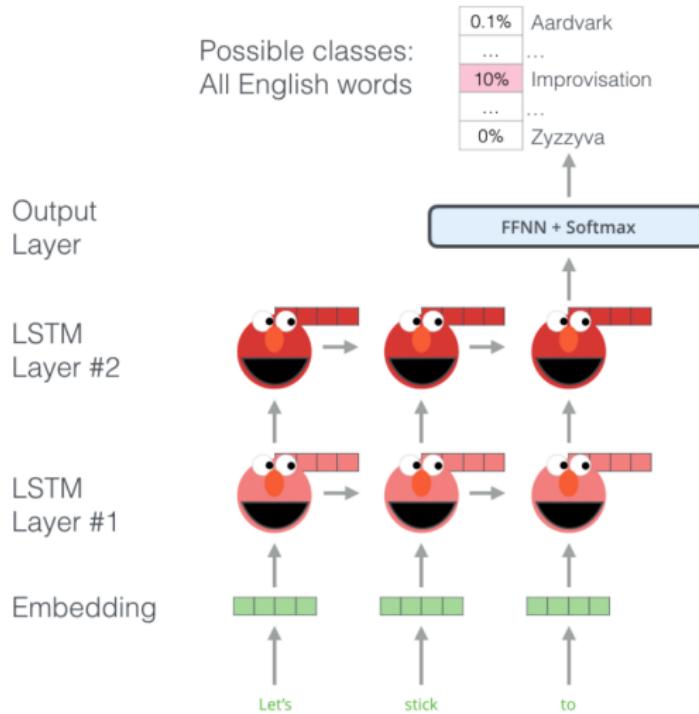
# CoVe (McCann et al., 2017)

- Key idea: Train a standard neural machine translation model
- Take the encoder directly as contextualized word embeddings
- Problems:
  - Translation requires paired (labeled) data
  - The embeddings are tailored to particular translation corpuses



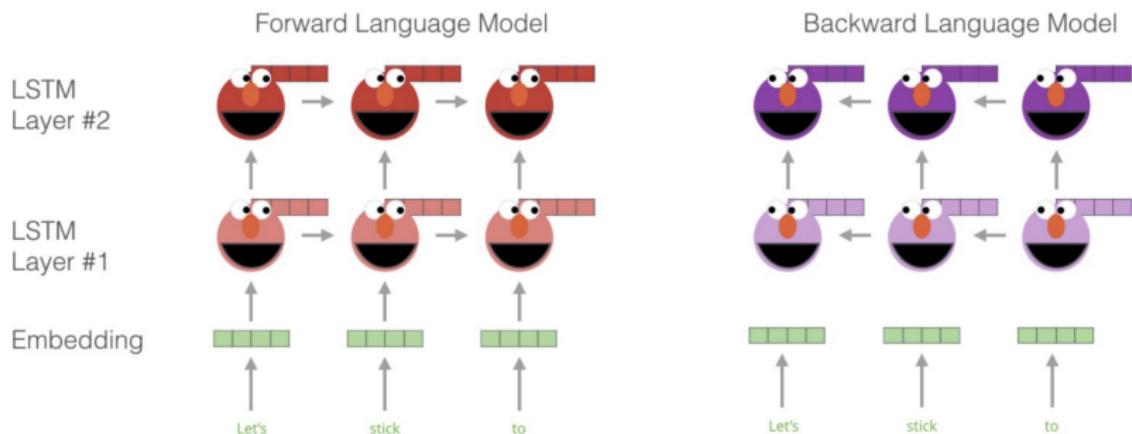
# Language model pretraining task

- Predict the next word given the prefix
- Can be defined on any unlabeled document



# ELMo (Peter et al., 2018)

- Key ideas:
  - Train a forward and backward LSTM language model on large corpus
  - Use the hidden states for each token to compute a vector representation of each word
  - Replace the word embedding by Elmo's embedding (with fixed Elmo's LSTM weights)



# ELMo results

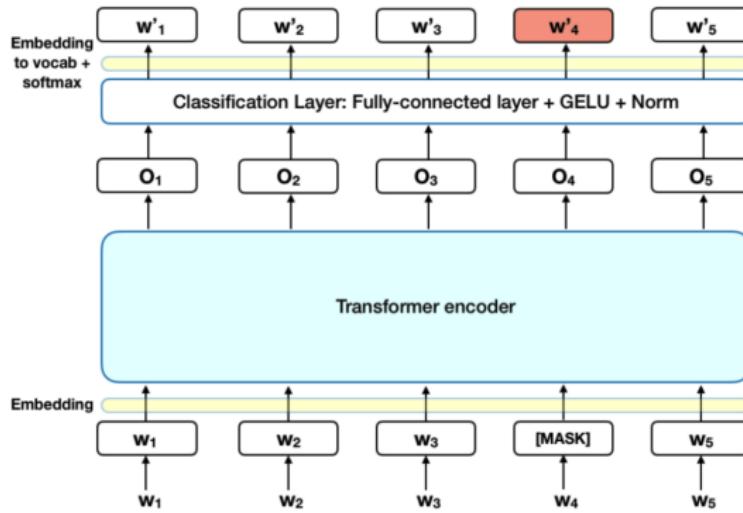
TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	$88.7 \pm 0.17$	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	$91.93 \pm 0.19$	90.15	$92.22 \pm 0.10$	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	$54.7 \pm 0.5$	3.3 / 6.8%

# BERT

- Key idea: replace LSTM by Transformer
- Define the generated pretraining task by masked language model
- Two pretraining tasks
- Finetune both BERT weights and task-dependent model weights for each task

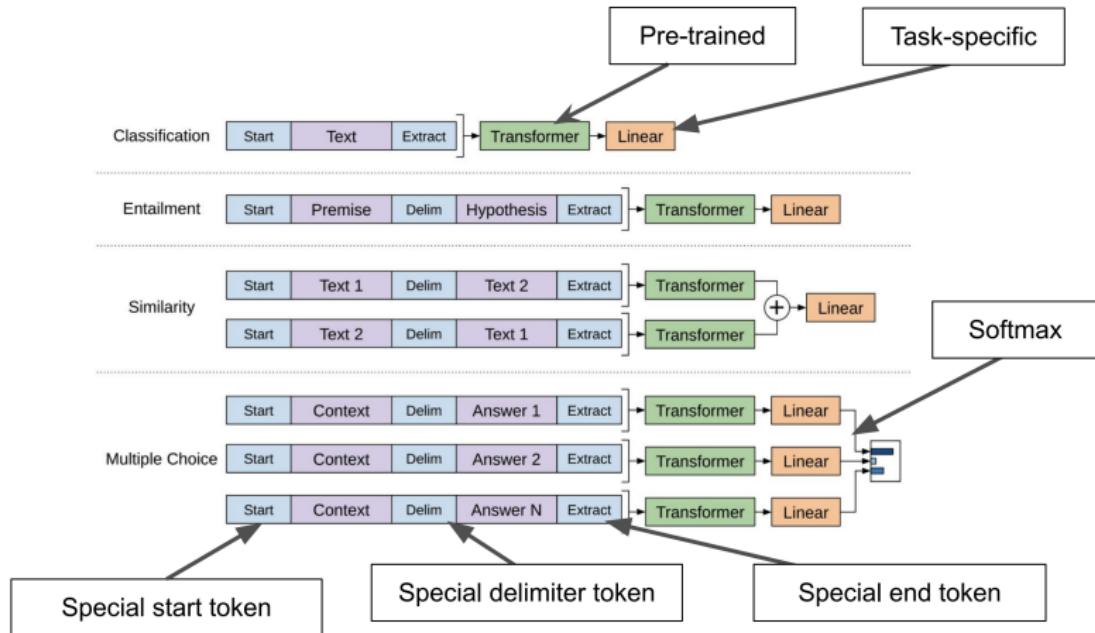
# BERT pretraining loss

- Masked language model: predicting each word by the rest of sentence
- Next sentence prediction: the model receives pairs of sentences as input and learns to predict if the second sentence is the subsequent sentence in the original document.



# BERT finetuning

- Keep the pretrained Transformers
- Replace or append a layer for the final task
- Train the whole model based on the task-dependent loss



# BERT results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

- BERT base: 110M parameters, BERT large: 340M parameters